

Suunnitteludokumentti

Boa Open Access

Helsinki 3.5.2006

Ohjelmistotuotantoprojekti

HELSINGIN YLIOPISTO

Tietojenkäsittelytieteen laitos

Kurssi

581260 Ohjelmistotuotantoprojekti (6 ov)

Projektiryhmä

Ilmari Heikkinen
Timo Hintsu
Erno Härkönen
Arto Vuori
Mikko Kautto

Asiakas

Olli Niinivaara

Johtoryhmä

Juha Taina
Riikka Kaven

Kotisivu

<http://www.cs.helsinki.fi/group/boa>

Versiohistoria

Versio	Päiväys	Tehdyt muutokset
1.0	20.3.2006	Ensimmäinen versio
2.0	3.4.2006	Toisen iteraation aiheuttamia muutoksia
2.5	7.4.2006	Korjauksia ja lisäyksiä
2.7	11.4.2006	Ulkoasua hiottu
2.8	14.4.2006	Korjauksia
2.9	3.5.2006	Korjauksia

Sisältö

1	Johdanto	1
2	Sanasto	1
3	Yleisarkkitehtuuri	1
4	Pääkomponentit	3
4.1	Tiedostonlukija	3
4.1.1	Luokat ja rajapintaluokat	3
4.1.2	Metodien paluuarvot	3
4.2	Citeseer-transformaattori	3
4.2.1	Luokat ja rajapintaluokat	4
4.2.2	Metodien paluuarvot	4
4.3	DBLP-transformaattori	4
4.3.1	Luokat ja rajapintaluokat	4
4.3.2	Metodien paluuarvot	5
4.4	Qs-transformaattori	5
4.4.1	Luokat ja rajapintaluokat	5
4.4.2	Metodien paluuarvot	5
4.5	QStatements-komponentti	5
4.5.1	Luokat ja rajapintaluokat	5
4.5.2	Metodien paluuarvot	7
4.5.3	QConnectionHandler	11
4.5.4	QStatementIndex	12
4.6	Statement-tulostaja	12
4.6.1	Luokat ja rajapintaluokat	12
4.6.2	Metodien paluuarvot	13
5	Apukomponentit	13
5.1	Metodien paluuarvot	15
6	Ulkoiset komponentit	16
7	Rajapinnat ja luokkien keskinäiset suhteet	16

8	Datalähde- ja käyttäjäkonfiguraatitiedostot	16
9	Pääohjelma	16
	Lähteet	16

1 Johdanto

Tämä on Boa Open Access -ryhmän kevään 2006 ohjelmistotuotantoprojektikurssilla toteuttaman Open Access -ohjelmiston suunnitteludokumentti. Suunnitteludokumentissa kuvataan ohjelmiston yleisarkkitehtuuri, osajärjestelmien väliset rajapinnat ja osajärjestelmät tilanteen vaatimalla tarkkuudella. Osajärjestelmien sisäistä toimintaa tietorakenteineen ja algoritmeineen on selitetty tarkemmin, kun se on ollut hyödyllistä. Suunnitteludokumenttia on kirjoitettu iteratiivisesti, koska projekti toteutetaan iteratiivisella tuotantomenetelmällä.

2 Sanasto

Qriterium Metadataympäristö [6]

Citeseer Open Access -tietokanta [2]

DBLP Digital Bibliography and Library Project. Tietojenkäsittelytiedettä käsittelevien artikkeleiden tietokanta. Huhtikuussa 2006 DBLP listasi yli 740 000 artikkelia. DBLP on perustettu 1980-luvulla. [3]

Transformaattori Muuntaja, muuttaja. Transformaattorin tehtävä on yhtenäistää erilaisia organisoitua metadataa yhtenäiseen muotoon.

Qs-lause Qriterium Statement. Qriterium-ympäristön metadatankuvausmenetelmä, johon erilaisilla organisoitu metadata muunnetaan. XML-muotoista.

QStatement Kuvaa yhtä Qs-lausetta. QStatement voi kuvata esimerkiksi dokumenttia, henkilöä, organisaatiota tai tietolähdettä.

3 Yleisarkkitehtuuri

Ohjelma koostuu kuudesta pääkomponentista ja ulkoisista komponenteista, joista kerrotaan luvussa "Ulkoiset komponentit".

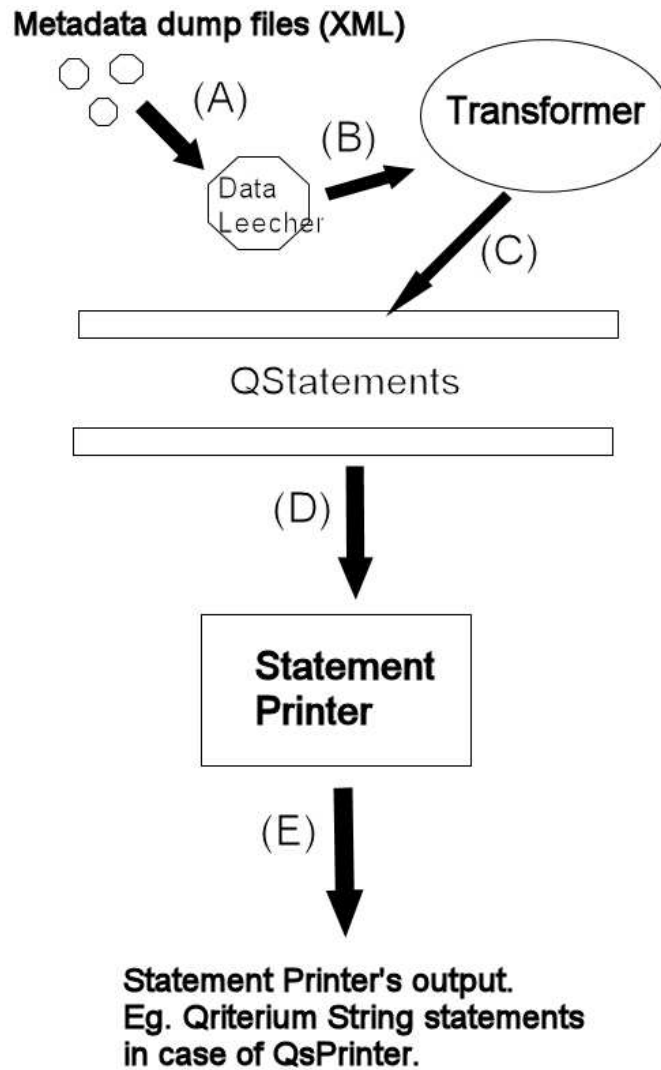
Tiedostonlukija Lukee levytä dataa ja tarjoaa sen eteenpäin.

Citeseer-transformaattori Muuntaa Citeseerin metadataa QStatementeiksi.

DBLP-transformaattori Muuntaa DBLP:n metadataa QStatementeiksi.

Qs-transformaattori Lukee Qs-lauseita järjestelmään. Muuntaa lauseet takaisin QStatementeiksi.

QStatements-komponentti Vastaanottaa transformaattorilta QStatementteja ja pitää niiden keskinäiset yhteydet ehjinä. Luovuttaa QStatementit edelleen Statement-tulostajille.



Kuva 1: Ohjelman yleisarkkitehtuuri

Statement-tulostaja Tulostaa QStatementteja. Esimerkiksi Qs-tulostaja, joka on yksi Statement-tulostajan erikoistapaus, tulostaa Qs-lauseita joko pakkaamattomana tai gzip-pakattuna. Se voi pilkkoa tulosteen halutun kokoisiksi palasiksi.

Kuvassa 1 nuolet kuvaavat datan kulkusuuntaa. DataLeecher lukee metadataa jostain lähteestä, esimerkiksi levytä, (A) ja luovuttaa datan transformaattorille InputStreamina (B). Transformaattori luo QStatementteja saadusta datasta ja antaa niitä QStatements-komponentille (C). QStatements-komponenttia kuuntelevat statement-tulostajat saavat tiedon uusista QStatementeista reaaliajassa (D) ja tulostavat QStatementit oman toteutuksensa mukaisesti. Esimerksiksi Qs-tulostajan tulosteena saadaan (E) Qs-lauseita (XML-tiedostoja).

4 Pääkomponentit

Tässä osassa on pyritty valottamaan komponenttien tehtäviä tarkemmin. Komponentit on jaettu luokkiin ja rajapintaluokkiin (Javan "interface"), jotka esitellään tässä osassa. Luokkien rajapinnat kuvataan tarkemmin seuraavassa luvussa.

4.1 Tiedostonlukija

Tässä ohjelmassa metadata-tiedostot haetaan levyiltä, eikä esimerkiksi OAI-PMH-protokollalla suoraan internetistä. Tiedostonlukija-komponentti huolehtii datan saamisesta eteenpäin transformaattoreille.

Metadataa on saatavilla mitä erilaisimmissa formaateissa ja sitä voidaan saattaa transformaattoreille erilaisilla tavoilla erilaisten protokollien yli. Tiedostonlukija kapsuloi datan hankintakeinot ja tarjoaa datan eteenpäin `InputStream`ina. Transformaattoreiden näkökulmasta jostain tulee vain datavirtaa, eikä niiden tarvitse erotella sitä, onko data alunperin tullut levyiltä tiedostosta vai jollain muulla tavalla.

4.1.1 Luokat ja rajapintaluokat

DataLeecher Rajapinta, jonka tiedostonlukija toteuttaa. Toteuttajat kapsuloivat sen, miten ja mistä data on hankittu, ja tarjoavat transformaattorille `InputStream`in.

FileReader Tiedostonlukija-luokka, joka toteuttaa `DataLeecher`-rajapinnan.

4.1.2 Metodien paluuarvot

Luokka	Metodi	paluuarvon selite
<code>DataLeecher</code>	<code>abstract getInputStream()</code>	<code>InputStream</code>
<code>DataLeecher</code>	<code>abstract getFileName()</code>	<code>String</code> , luettavan tiedoston nimi
<code>FileReader</code>	<code>getInputStream()</code>	<code>InputStream</code>
<code>FileReader</code>	<code>getFileName()</code>	<code>String</code>

4.2 Citeseer-transformaattori

Citeseer-transformaattori erottelee saamastaan XML-datasta tietueet, jotka on merkitty citeseer-datassa record-tagein. Olemme huomanneet, että Citeseerin tarjoamassa datassa kaikki recordit eivät ole validia XML:ää. Vialliset recordit hylätään. Kun viallinen record osuu kohdalle, siirrytään seuraavaan.

XML:n jäsentämisessä käytetään XPathia (DOM-pohjainen jäsentäjä).

Transformaattoreiden pitää pystyä erottelemaan syötteestään kiinnostava data. Qriterium Statement (ks. Vaatimusdokumentti [1]) määrää, mikä on kiinnostavaa dataa.

Transformaattori luo oikean tyyppisiä QStatement-instansseja jäsentämistään recordeista. Oikealla tyyppillä tarkoitetaan sitä, että transformaattorin täytyy tunnistaa, mitä recordissa kuvataan. Kuvauksen kohteena voi olla itse dokumentin lisäksi muuttuva määrä tekijöitä. Tekijät voivat taas olla ihmisiä tai organisaatioita. Dokumenttien välillä on yhteyksiä. Yhteydet voivat olla esimerkiksi dokumenttien välisiä. Myös muunlaisia yhteyksiä on. Dokumentin ja sen tekijän välillä on aina yhteys. Transformaattorin on kyettävä tunnistamaan, minkä tyyppinen yhteys on kyseessä ja luoda vastaavan tyyppisiä QStatement-instansseja.

4.2.1 Luokat ja rajapintaluokat

Transformer Abstrakti yläluokka kaikille transformaattoreille

CiteseerTransformer Citeseer-metadataa ymmärtävä transformaattori.

4.2.2 Metodien paluarvot

Luokka	Metodi	paluarvon selite
Transformer	transform(DataLeecher, QStatements)	<i>void</i> , suorittaa muuntamisen
Transformer	getValidItems()	<i>long</i> , ehjien muunnosten määrä
Transformer	getInvalidItems()	<i>long</i> , epäonnistuneiden muunnosten määrä
CiteseerTransformer	transform(DataLeecher, QStatements)	<i>void</i> , suorittaa muuntamisen
CiteseerTransformer	getValidItems()	<i>long</i> , ehjien muunnosten määrä
CiteseerTransformer	getInvalidItems()	<i>long</i> , epäonnistuneiden muunnosten määrä

4.3 DBLP-transformaattori

DBLP-transformaattori tuo DBLP:n tarjoamia artikkeleita järjestelmään. Transformaattori toteutetaan hyvin samalla tavalla kuin Citeseer-transformaattori. Sisään luettu data on eri muodossa kuin Citeseer-data, mutta tämänkin transformaattorin tarkoituksena on luoda syötteestään QStatement-instansseja.

DBLP-transformaattori on yleisen Transformer-luokan erikoistapaus.

4.3.1 Luokat ja rajapintaluokat

DBLPTransformer DBLP-metadataa ymmärtävä transformaattori.

4.3.2 Metodien paluuarvot

Luokka	Metodi	paluuarvon selite
DBLPTransformer	transform(DataLeecher, QStatements)	<i>void</i> , suorittaa muuntamisen
DBLPTransformer	getValidItems()	<i>long</i> , ehjien muunnosten määrä
DBLPTransformer	getInvalidItems()	<i>long</i> , epäonnistuneiden muunnosten määrä

4.4 Qs-transformaattori

Qs-transformaattori osaa muuntaa Qs-tulostajan tuottamia Qs-lauseita takaisin QState-menteiksi. Qs-transformaattori on yleisen Transformer-luokan erikoistapaus.

4.4.1 Luokat ja rajapintaluokat

QsTransformer Qs-lauseita ymmärtävä transformaattori.

4.4.2 Metodien paluuarvot

Luokka	Metodi	paluuarvon selite
QsTransformer	transform(DataLeecher, QStatements)	<i>void</i> , suorittaa muuntamisen
QsTransformer	getValidItems()	<i>long</i> , ehjien muunnosten määrä
QsTransformer	getInvalidItems()	<i>long</i> , epäonnistuneiden muunnosten määrä

4.5 QStatements-komponentti

QStatements vastaanottaa transformaattorilta QStatementteja ja pitää niiden keskinäiset yhteydet ehjinä. QStatements-komponenttia kuuntelevat statement-tulostajat saavat QStatements-komponentilta tiedon uusista QStatementeista reaaliajassa.

4.5.1 Luokat ja rajapintaluokat

QStatements Vastaanottaa QStatementteja ja ilmoittaa niistä kuunteleville statement-tulostajille.

QStatement Kuvaa yhtä Qriterium-statementtia.

QHeader Kuvaa Qriterium-statementin otsaketta.

QDocument QStatementin ilmentymä. Kuvaa dokumentin metadattaa.

QActor QStatementin ilmentymä. Kuvaa henkilöä tai organisaatiota. Jakautuu alatyyppeihin Person ja Organization.

Person Kuvaa henkilöä.

Organization Kuvaa organisaatiota.

QConnection QStatementin ilmentymä. Kuvaa yhteyttä kahden Qriterium-statementin välillä. QConnectionit on edelleen jaettu alatyyppeihin (CreatedForConnection, ContactedAtConnection, PublishedByConnection, StoredByConnection, SupportedByConnection, CreatedInConnection, ManagedByConnection, IdentifiedAsConnection, ConsistsOfConnection, AwareOfConnection, CreatedByConnection).

QSource QStatementin ilmentymä. Kuvaa "raakadatan"lähdettä. Raakadataa on esimerkiksi Citeseerin metadatadumpit.

QContent Qstatementin ilmentymä. Kuvaa metadataa dokumentin sisällöstä.

QStatementIndex Indeksoi Boan sisäiset tunnisteet käyttäen avaimena raa'an metadata-syötteen omia tunnisteita. QStatementIndexiltä voidaan myöhemmin selvittää, onko joku QStatement jo ennestään tuttu järjestelmälle.

QConnectionHandler Luo QConnectioneja sitä mukaan, kun se on mahdollista. QConnectionia ei voi luoda kahden QStatementin välille, ennen kuin kummatkin osapuolet on luettu sisään ohjelmaan. Säilyttää tilansa eräajojen välillä toteuttamansa Storable-rajapinnan avulla.

RawConnection QConnectionHandlerin tarvitsema apuluokka. Transformaattorit eivät luo suoraan QConnectioneja vaan RawConnectioneja, koska Transformaattori ei tiedä, tuleeeko yhteyden toista osapuolta löytymään ikinä tulevaisuudessa. QConnectionHandlerin vastuulle jää oikeiden QConnectionien luominen saamiensa RawConnection-olioiden ja QStatements-komponentin antamien signaalien perusteella.

QConnectionFactory Luo uusia halutun tyyppisiä QConnectioneja tyyppimerkkijonon perusteella.

IdFactory Huolehtii QStatementtejen sisäisten tunnisteiden yksilöllisyydestä. Luovuttaa pyydetessä uniikin id:n. Säilyttää tilansa eräajojen välillä toteuttamansa Storable-rajapinnan avulla.

XMLDataElement QStatementin sisäiset elementit, eli avain-arvo-parit on saatavissa ulos XMLDataElement-muodossa.

4.5.2 Metodien paluuarvot

Luokka	Metodi	paluuarvon selite
QStatements	addQStatement(QStatement)	<i>boolean</i> , true jos QStatement välitettiin statement-tulostajille, muutoin false.
QStatements	addQStatement(RawConnection)	<i>boolean</i> , true jos QStatement välitettiin statement-tulostajille.
QStatements	getConnectionHandler()	<i>QConnectionHandler</i> , testausta varten.
QStatements	getStatementIdByRawId(String)	<i>Long</i> null, ohjelman sisäinen QStatement-tunniste Long-oliona, jos indeksi löytää QStatementin annetulla rawId:llä. Muutoin palautetaan null.
QStatements	getIndex()	<i>QStatementIndex</i> , testausta varten.
QStatements	addPrinter(StatementPrinter)	<i>void</i> , lisää statement-tulostajan QStatements-komponentin kuuntelijaksi.
QStatement	setHeader(QHeader)	<i>void</i> , asettaa headerin
QStatement	getHeader()	<i>QHeader</i> , palauttaa headerin
QStatement	getId()	<i>long</i> , QStatementin sisäinen tunniste
QStatement	getHash()	<i>String</i> , alkuperäisen datalähteen tunniste-arvo
QStatement	getRawId()	<i>String</i> , alkuperäisen datalähteen antama tunniste
QStatement	getAt()	<i>Date</i> , QStatementin luontipäivämäärä
QStatement	isConnectable()	<i>boolean</i> , true jos QStatement on yhdistettävissä QConnectionilla, muutoin false.
QStatement	getXMLNamespace()	<i>String</i> , XML-nimiavaruus
QStatement	toString()	<i>String</i> , merkkijonoesitysmuoto
QStatement	abstract getXMLDataElements()	<i>XMLDataElement[]</i> , XML-esitysmuodon avain-arvoparit
QStatement	abstract getStartTag()	<i>String</i> , XML-esityksen vaatima aloitustagi
QStatement	abstract getEndTag()	<i>String</i> , XML-esityksen vaatima lopetus-tagi
QHeader	getId()	<i>long</i> , QStatementin sisäinen tunniste
QHeader	getHash()	<i>String</i> , alkuperäisen datalähteen tunniste-arvo
QHeader	getRawId()	<i>String</i> , alkuperäisen datalähteen antama tunniste
QHeader	getAt()	<i>Date</i> , QStatementin luontipäivämäärä
QHeader	getXMLDataElements()	<i>XMLDataElement[]</i> , XML-esitysmuodon avain-arvoparit

Luokka	Metodi	paluuarvon selite
QDocument	setTitle(String)	<i>void</i>
QDocument	setLanguage(String)	<i>void</i>
QDocument	setFilename(String)	<i>void</i>
QDocument	getTitle()	<i>String</i> , dokumentin otsikko
QDocument	getLanguage()	<i>String</i> , dokumentin kieli
QDocument	getFilename()	<i>String</i> , dokumentin sijainti (URI)
QDocument	getStartTag()	<i>String</i> , XML-esityksen vaatima aloitustagi
QDocument	getEndTag()	<i>String</i> , XML-esityksen vaatima lopetus- tagi
QDocument	getXMLDataElements()	<i>XMLDataElement[]</i> , XML- esitysmuodon avain-arvoparit
QActor	setName(String)	<i>void</i>
QActor	getName()	<i>String</i> , actorin nimi
QActor	isConnectable()	<i>boolean</i> , aina true, koska QActorit on aina yhdistettävissä QConnectionilla
QActor	getStartTag()	<i>String</i> , XML-esityksen vaatima aloitustagi
QActor	getEndTag()	<i>String</i> , XML-esityksen vaatima lopetus- tagi
QActor	getXMLDataElements()	<i>XMLDataElement[]</i> , XML- esitysmuodon avain-arvoparit
QActor	abstract getSubType-String()	<i>String</i> , actorin alatyyppi
Person	getSubTypeString()	<i>String</i> , henkilön oma alatyypitunniste
Organization	getSubTypeString()	<i>String</i> , organisaation oma alatyypitunniste
QConnection	setFrom(long)	<i>void</i>
QConnection	setFrom(Long)	<i>void</i>
QConnection	setTo(long)	<i>void</i>
QConnection	setTo(Long)	<i>void</i>
QConnection	getFrom()	<i>Long</i> , yhteyden viittaavan pään sisäinen tunniste
QConnection	getTo()	<i>Long</i> , yhteyden viitattavan pään sisäinen tunniste
QConnection	getStartTag()	<i>String</i> , XML-esityksen vaatima aloitustagi
QConnection	getEndTag()	<i>String</i> , XML-esityksen vaatima lopetus- tagi
QConnection	getXMLDataElements()	<i>XMLDataElement[]</i> , XML- esitysmuodon avain-arvoparit
QConnection	toString()	<i>String</i> , merkkijonoesitysmuoto yhteydestä
QConnection	abstract getType-String()	<i>String</i> , yhteyden tunnistemerkkijono

Luokka	Metodi	paluarvon selite
AwareOfConnection	getTypeString()	<i>String</i> , yhteyden tunnistemerkkijono
ConsistOfConnection	getTypeString()	<i>String</i> , yhteyden tunnistemerkkijono
ContactedAtConnection	getTypeString()	<i>String</i> , yhteyden tunnistemerkkijono
CreatedByConnection	getTypeString()	<i>String</i> , yhteyden tunnistemerkkijono
CreatedForConnection	getTypeString()	<i>String</i> , yhteyden tunnistemerkkijono
CreatedInConnection	getTypeString()	<i>String</i> , yhteyden tunnistemerkkijono
IdentifiedAsConnection	getTypeString()	<i>String</i> , yhteyden tunnistemerkkijono
ManagedByConnection	getTypeString()	<i>String</i> , yhteyden tunnistemerkkijono
PublishedByConnection	getTypeString()	<i>String</i> , yhteyden tunnistemerkkijono
StoredByConnection	getTypeString()	<i>String</i> , yhteyden tunnistemerkkijono
SupportedByConnection	getTypeString()	<i>String</i> , yhteyden tunnistemerkkijono
QSource	setName(String)	<i>void</i>
QSource	setOpenURL(URL)	<i>void</i>
QSource	getName()	<i>String</i>
QSource	getOpenURL()	<i>URL</i>
QSource	setEncoding(String)	<i>void</i>
QSource	getEncoding()	<i>String</i>
QSource	setFormat(String)	<i>void</i>
QSource	getFormat()	<i>String</i>
QSource	setTerminology(String)	<i>void</i>
QSource	getTerminology()	<i>String</i>
QSource	setCompression(String)	<i>void</i>
QSource	getCompression()	<i>String</i>
QSource	setBackprotocol(String)	<i>void</i>
QSource	getBackProtocol()	<i>String</i>
QSource	setTopprotocol(String)	<i>void</i>
QSource	getTopprotocol()	<i>String</i>
QSource	setLaunched(Date)	<i>void</i>
QSource	getLaunched()	<i>Date</i>
QSource	getStartTag()	<i>String</i> , XML-esityksen vaatima aloitustagi
QSource	getEndTag()	<i>String</i> , XML-esityksen vaatima lopetus- tagi
QSource	getXMLDataElements()	<i>XMLDataElement[]</i> , XML- esitysmuodon avain-arvoparit
QContent	isConnectable()	<i>boolean</i> , aina true, koska QContent on yhdistettävissä QConnectionilla
QContent	toString()	<i>String</i> , merkkijonoesitysmuoto QContentista
QContent	getStartTag()	<i>String</i> , XML-esityksen vaatima aloitustagi
QContent	getEndTag()	<i>String</i> , XML-esityksen vaatima lopetus- tagi
QContent	getXMLDataElements()	<i>XMLDataElement[]</i> , XML- esitysmuodon avain-arvoparit

Luokka	Metodi	paluarvon selite
QStatementIndex	put(QStatement)	<i>void</i> , lisää QStatementin indeksiin
QStatementIndex	putRaw(String s, Long l)	<i>void</i> , lisää indeksiin alkion, jonka avain on s ja arvo on l
QStatementIndex	get(String s)	<i>Long</i> <i>null</i> , palauttaa Long-arvon, jos avaimella s on alkio indeksissä. Jos mitään ei löydy, palauttaa null
QStatementIndex	store()	<i>boolean</i> , true jos QStatementIndex'n tilan tallennus onnistui, muutoin false.
QStatementIndex	store(File)	<i>boolean</i> , sama kuin yllä, mutta tallennustiedosto on määriteltävissä. Metodi on testejä varten tehty.
QStatementIndex	retrieve()	<i>boolean</i> , true jos QStatementIndex'n tilan palautus onnistui, muutoin false.
QStatementIndex	retrieve(File)	<i>boolean</i> , sama kuin yllä, mutta palautustiedosto on määriteltävissä. Metodi on testejä varten tehty.
QConnectionHandler	addRawConnection(RawConnection)	<i>void</i> , lisää uuden RawConnectionin QConnectionHandleriin
QConnectionHandler	statementAdded(String)	<i>int</i> , luotujen QConnection-olioiden määrä
QConnectionHandler	store()	<i>boolean</i> , true jos QConnectionHandler'n tilan tallennus onnistui, muutoin false.
QConnectionHandler	store(File)	<i>boolean</i> , sama kuin yllä, mutta tallennustiedosto on määriteltävissä. Metodi on testejä varten tehty.
QConnectionHandler	retrieve()	<i>boolean</i> , true jos QConnectionHandler'n tilan palautus onnistui, muutoin false.
QConnectionHandler	retrieve(File)	<i>boolean</i> , sama kuin yllä, mutta palautustiedosto on määriteltävissä. Metodi on testejä varten tehty.
QConnectionHandler	getHashSize()	<i>int</i> , RawConnection-olioiden määrä QConnectionHandlerissä
QConnectionHandler	getRawConnections()	<i>RawConnection[]</i> , palauttaa RawConnectionit taulukossa. Tämä metodi on testejä varten olemassa.
RawConnection	getTo()	<i>String</i> , viittaavan statementin rawId
RawConnection	getFrom()	<i>String</i> , viitattavan statementin rawId
RawConnection	getType()	<i>String</i> , yhteyden tyypin merkkijonoesitysmuoto
RawConnection	toString()	<i>String</i> , merkkijonoesitysmuoto
RawConnection	getStartTag()	<i>String</i> , XML-esityksen vaatima aloitustaggi
RawConnection	getEndTag()	<i>String</i> , XML-esityksen vaatima lopetus-taggi
RawConnection	getXMLDataElements()	<i>XMLDataElement[]</i> , XML-esitysmuodon avain-arvoparit
QConnectionFactory	getInstance(String type)	<i>QConnection</i> <i>null</i> , jos type on validi QConnection-tyyppi palauttaa pyydetyn tyyppisen QConnection-olion. Jos tyyppi on tuntematon palautetaan null.

Luokka	Metodi	paluuarvon selite
IdFactory	getUniqueId	<i>long</i> , uniikki tunniste QStatement-oliota varten
IdFactory	store()	<i>boolean</i> , true jos IdFactory'n tilan tallennus onnistui, muutoin false.
IdFactory	retrieve()	<i>boolean</i> , true jos IdFactory'n tilan palautus onnistui, muutoin false.
XMLDataElement	getNamespace()	<i>String</i> , XML-elementin nimiavaruus
XMLDataElement	getName()	<i>String</i> , XML-elementin nimi
XMLDataElement	getValue()	<i>String</i> , XML-elementin arvo

4.5.3 QConnectionHandler

QStatements-komponentin QConnectionHandlerin toimintaa on syytä selventää lisää.

Tehtävät

- Luo QConnection-instansseja.
- Ylläpitää vaillinaisia QConnectioneja (RawConnectioneja)

Toteutus QConnectionHandler:lle annetaan RawConnection-olioita, jotka sisältävät kahden QDocumentin rawId:t. Handler ylläpitää RawConnection-olioita TreeMap-tietorakenteessa (Sunin toteutus punamustasta puusta). TreeMapiin on päädytty, koska RawConnectioneja pitää pystyä jatkuvasti etsimään, lisäämään ja poistamaan tietorakenteesta tehokkaasti. Punamusta puu takaa takaa muokkausoperaatioille ajan $O(n \log n)$ ja etsimiselle aikaluokan $O(\log n)$, jossa n on solmujen määrä puussa. Puun avaimena on RawConnectionista saatu toRawId (String). QStatements-olio antaa uuden RawConnectionin Handlerille aina kun transformaattori sellaisen sille luovuttaa. Transformaattorin pitää antaa RawConnection-instanssille QHeader, josta myöhemmin kaivetaan osa QConnectionin otsaketiedoista, ja QConnectionin tyyppi (String).

Kun transformaattori löytää syötteestään uuden statementin, tieto siitä tulee QStatements-olion kautta Handlerille metodilla statementAdded(QStatement statement). Handleri etsii RawConnection-olioita puurakenteestaan käyttäen syötetyn QStatementin rawId:tä. Jos RawConnection löytyy, tarkoittaa se sitä, että voidaan luoda uusi QConnection-instanssi. Uusi QConnection saa headertiedoistaan RawConnectionilta seuraavat kentät: qs:rawId, qs:origin, qs:at. Handleri luo oikeantyyppisen QConnection-instanssin QConnectionFactoryin avulla käyttäen RawConnectioniin aikaisemmin säilöttyä tyyppi-stringiä.

Lopulta QConnectionHandler luovuttaa uuden QConnectionin QStatements-komponentille, joka ilmoittaa normaalisti uudesta QStatementista kuunteleville statement-tulostajille. Tätä tehdään, kunnes QConnectionHandlerin RawConnection-puusta ei enää löydy solmuja ko. toRawId:llä.

4.5.4 QStatementIndex

QStatementIndexin tehtävä on pitää kirjaa QStatementeista (sellaisista, joihin on mahdollista liittää yhteys "QConnection"), jotta on mahdollista sekä eräajon aikana että seuraavilla ajokerroilla selvittää, onko jokin uudelleen vastaan tullut QStatement jo ennestään tuttu järjestelmälle, eli onko sillä jo sisäinen tunniste. Mikäli sisäinen tunniste löytyy jo, ei tarvitse luoda uutta QStatementtia vaan käytetään jo ennen luotua.

QStatementtien uutuus selvitetään kysymällä QStatementIndexiltä; löytyykö tietylle rawId:lle jo merkintää. QStatementIndex palauttaa vastauksena sisäisen tunnisteiden, tai null jos vastaavuutta ei löydy.

QStatements-olio ilmoittaa QStatementIndexille uusista QStatementeista kutsumalla indeksin put(String rawId, QStatement statement) -metodia. Koska indeksin pitää säilyttää tilansa eräajojen välissä, sen tila pitää olla tallennettavissa levyllä. QStatementIndex toteuttaa rajapintaluokan Storable, jonka metodit store() ja retrieve() mahdollistavat halutun toiminnallisuuden. Talletusmuotona käytetään XML:ää.

4.6 Statement-tulostaja

Tämän ohjelman tarkoitus on pääasiassa tuottaa Qriterium String -muotoista XML-ulostuloa - eli Qs-lauseita. Tässä kappaleessa selvitetäänkin Qs-tulostajan toimintaa.

Qs-tulostajan tehtävä on luoda Qriterium-muotoista XML-dataa (ks. vaatimusdokumentti [1]) QStatements-komponentilta saaduista QStatement-olioista. Qs-tulostaja luo XML-tiedostoja halutun kokoisissa palasissa. Myös gzip-pakattu ulostulo on mahdollista.

4.6.1 Luokat ja rajapintaluokat

StatementPrinter Rajapinta statement-tulostajille

QSPrinter Tulostaa Qriterium-muotoisia XML-tiedostoja QStatement-olioista.

XMLProperties Luokka, joka tarjoaa Qriterium-XML:n tarvitsemia vakioita.

4.6.2 Metodien paluuarvot

Luokka	Metodi	paluuarvon selite
StatementPrinter	abstract print(QStatement)	<i>void</i> , tulostaa parametrinä annetun QStatementin
StatementPrinter	abstract close()	<i>void</i> , viimeistelee tulostuksen
QsPrinter	print(QStatement)	<i>void</i> , tulostaa QStatementin Qs-muodossa
QsPrinter	close()	<i>void</i> , viimeistelee tulostuksen
QsPrinter	getFilePath()	<i>String</i> , tulostetiedostojen hakemisto
QsPrinter	setFilePath(String)	<i>void</i> , asettaa tulostushakemiston
QsPrinter	getFileName()	<i>String</i> , tulostiedostojen nimi (etuliite)
QsPrinter	setFileName(String)	<i>void</i> , asettaa tulostiedoston nimen (etuliiteen)
QsPrinter	isGzip()	<i>boolean</i> , true jos tulostetaan gzip-pakattua dataa, muuten false
QsPrinter	setGzip(boolean)	<i>void</i>
QsPrinter	getOutputSize()	<i>long</i> , tulostustiedostojen maksimikoko tavuina
QsPrinter	setOutputSize(long)	<i>void</i> , asettaa tulostustiedostojen maksimikoon tavuina
QsPrinter	getDataWriter()	<i>DataWriter</i> , (XML-DataWriter)
QsPrinter	setDataWriter(DataWriter)	<i>void</i> , asettaa XML-DataWriter'n
QsPrinter	getOutputStreamWriter()	<i>OutputStreamWriter</i>
QsPrinter	getCurrentFileName()	<i>String</i> , nykyinen tiedostonimi, johon tulostetaan

5 Apukomponentit

Apuluokat, jotka eivät varsinaisesti ole minkään muun komponentin osia, esitellään tässä kappaleessa.

LogFormatter Lokitietojen kirjaaja

StatusPrinter Edistymisilmaisimien käyttäjää varten

StringUtilities Apu-string-funktioita

Storable Rajapinta olion tilan säilyttämiseen ja palauttamiseen

Connection Rajapinta yksisuuntaisille yhteyksille

Source Metadatalähteestä

User Metadatalähteestä

SourceState Kuvaa metadatalähteen tilaa eräajojen välillä ja ylläpitää FileState-olioita

FileState Kuvaa yhden syötetiedoston tilaa eräajoen jälkeen

Loader Apuluokka konfiguraatitiedostojen lukemiseen

5.1 Metodien paluuarvot

Luokka	Metodi	paluuarvon selite
LogFormatter	format(LogRecord)	<i>String</i> , lokitieto
StatusPrinter	printTemp(String)	<i>void</i>
StatusPrinter	print(String)	<i>void</i>
StatusPrinter	println(String)	<i>void</i>
StatusPrinter	printProgressChar()	<i>void</i>
StringUtilities	zeroFill(long, int)	<i>String</i> , <long> merkkiä pitkä merkkijono mahdollisine etunollineen
Storable	store()	<i>boolean</i> , true jos olion tilan tallennus onnistui, muutoin false.
Storable	retrieve()	<i>boolean</i> , true jos olion tilan palautus onnistui, muutoin false.
Connection	getTo()	<i>Object</i> , yhteyden viittaaja
Connection	getFrom()	<i>Object</i> , yhteyden viitattava
Source	getFileName()	<i>String</i>
Source	getName()	<i>String</i>
Source	getUser()	<i>String</i>
Source	getInputDir()	<i>String</i>
Source	getOutputDir()	<i>String</i>
Source	getDataDir()	<i>String</i>
Source	getTransformer()	<i>String</i>
Source	getPrinter()	<i>String</i>
Source	getOutputFile()	<i>String</i>
Source	getGzip()	<i>String</i>
Source	getCutSize()	<i>long</i>
User	getName()	<i>String</i>
SourceState	getActiveState()	<i>SourceState</i>
SourceState	isCompleted()	<i>String</i>
SourceState	getFileNames()	<i>void</i>
SourceState	setFileState(String, FileState)	<i>void</i>
SourceState	getFileState(String)	<i>FileState</i>
SourceState	setNextFileID(String)	<i>void</i>
SourceState	setNextSerialID()	<i>void</i>
SourceState	setSourceID(long)	<i>void</i>
SourceState	getSourceID()	<i>long</i>
SourceState	getNextFileID()	<i>long</i>
SourceState	getNextSerialID()	<i>long</i>
SourceState	updateLastRun()	<i>void</i>
SourceState	getLastRun()	<i>Date</i>
SourceState	toString()	<i>String</i>
FileState	getCreationID()	<i>long</i>
FileState	setCreationID(long)	<i>void</i>
FileState	getDate()	<i>Date</i>
FileState	setDate(Date)	<i>void</i>
FileState	getInvalidItems()	<i>long</i>
FileState	setInvalidItems(long)	<i>void</i>
FileState	getValidItems()	<i>long</i>
FileState	setValidItems(long)	<i>void</i>
Loader	loadSettings(String)	<i>void</i>
Loader	loadSettings()	<i>void</i>

6 Ulkoiset komponentit

Ohjelmassa on hyödynnetty ulkoisia komponentteja tilanteissa, joissa juuri meidän tarkoitukseen sopiva komponentti on ollut vapaasti saatavilla.

XML-Writer XML-tulostaja. Hallitsee nimiavaruudet kätevästi. [4]

JArgs Komentoriviparametrijäsentäjä [5]

7 Rajapinnat ja luokkien keskinäiset suhteet

Yllä kuvattiin luokkien rajapinnat. Luokkien keskinäiset suhteet kuvataan erikseen UML-kaavion avulla. Kaaviota voidaan sanoa tarkaksi, sillä se kattaa miltei koko järjestelmän lukuun ottamatta muutamaa arkkitehtuurin kannalta epäolennaista luokkaa.

Kaavio tulee erikseen liitteenä (*boa_uml.pdf*).

8 Datalähde- ja käyttäjäkonfiguraatiodostot

Datalähteille - - esimerkiksi Citeseerille ja DBLP:lle - - luodaan konfiguraatiodostot, joissa säilytetään transformaattoreiden tarvitsemää tietoa. Konfiguraatiodostoissa pidetään kirjaa muunnettujen QStatementtien määrästä ko. lähteestä.

Ohjelman käyttäjille luodaan omat konfiguraatiodostot. Tiedosto sisältää käyttäjän yksilöivän tunnusteen ja muuta tietoa käyttäjästä, kuten henkilön nimen. Yhtä lähdekonfiguraatiota kohden voi olla yksi käyttäjäkonfiguraatio. Yksi käyttäjäkonfiguraatio voi liittyä moneen lähdekonfiguraatioon.

Konfiguraatiodostoja käytetään myös mainostarkoituksessa muunnetun Qs-datan yhteydessä. Ulkopuoliset saavat selville, kuka on tarjonnut tiedon.

9 Pääohjelma

Pääohjelman (org.qriterium.ui.Boa), tehtävänä on rakentaa Citeseer- DBLP- tai Qs-tyyppisestä datasta Qriterium String -muotoista XML-dataa (Qs-lauseita) käyttäen edellä mainittuja komponentteja.

Pääohjelma toimii konsolissa.

Lähteet

1 <http://www.cs.helsinki.fi/group/boa/docs/vaatimusdokumentti.pdf>

2 <http://citeseer.ist.psu.edu/>

3 <http://www.informatik.uni-trier.de/ley/db/>

4 <http://megginson.com/>

5 <http://jargs.sourceforge.net/>

6 <http://qriterium.org/>